



Managed by Fermi Research Alliance, LLC for the U.S. Department of Energy Office of Science

High Energy Physics Data Management for *Internet of Everything Topologies*

Federal Labs Technology Demonstrations and Discussions

Dr. Andrew J Norman

Fermilab Scientific Computing Division

My Research – Neutrino Oscillations

Experiments Designed to probe ***matter/anti-matter*** asymmetries

Send high power ν beams through Earth's crust

- Fermilab to Northern MN
- Measure transformations in beam

Designed to answer the next generation of fundamental particle physics (ν) questions

- Mass Hierarchy
- ν_3 dominant coupling (θ_{23} octant)
- CPV in ν sector
- Tests of 3-flavor mixing
- Detection of Supernovae ν 's



Research – Big Data

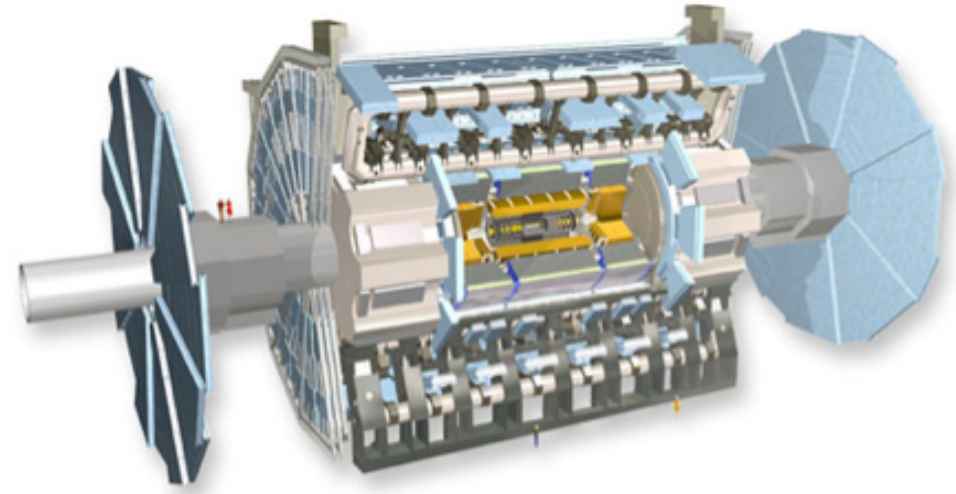
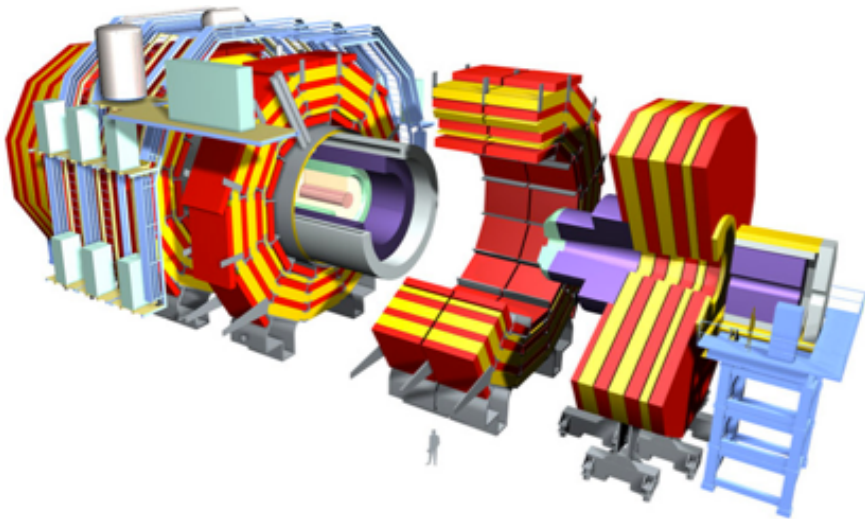
- *I will not tell you about matter, anti-matter, neutrinos, cosmic rays, supernova or magnetic monopoles*
- **Big Data** – Auxiliary research program developed to address large scale:
 - information management
 - data storage, retrieval and analysis
- Work designed to solve:
 - Peta+ scale data management
 - “stove pipe” problem in data sharing
 - Multi-experiment needs for information sharing and analysis integration
 - Rationalize & preserve legacy data sets

Research Team – Scientific Computing

- Data Management/Movement/Storage Research Groups
 - 12 full time researchers + ops staff
 - 10 physicists (PhD) [5 data management + 5 Storage]
 - 2 software engineers
 - 5 storage operations staff
- Focus on “Big Data” management and transport for analysis
 - Research is directed at satisfying experiment’s needs now and in the future (5 and 10 year horizons)
 - Coupled to research in distributed computing (grid & cloud infrastructures)
 - Develop the systems & tools
 - Integration with experiment’s infrastructure
 - Provide operational support for experiments

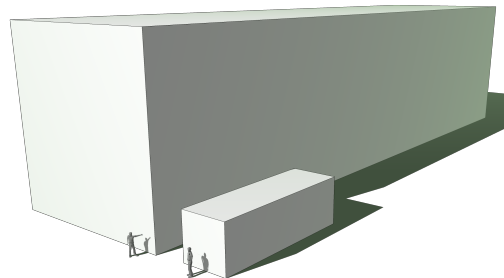
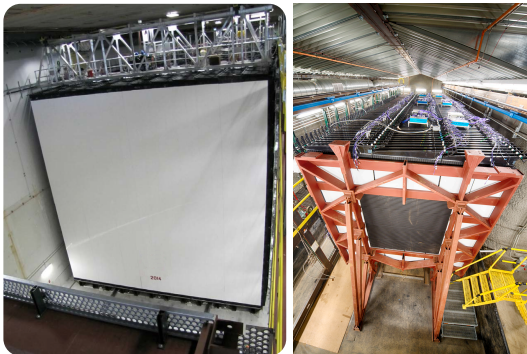
Data Volume Problem

CERN Collider



Atlas & CMS (2010-)

Fermilab Neutrino



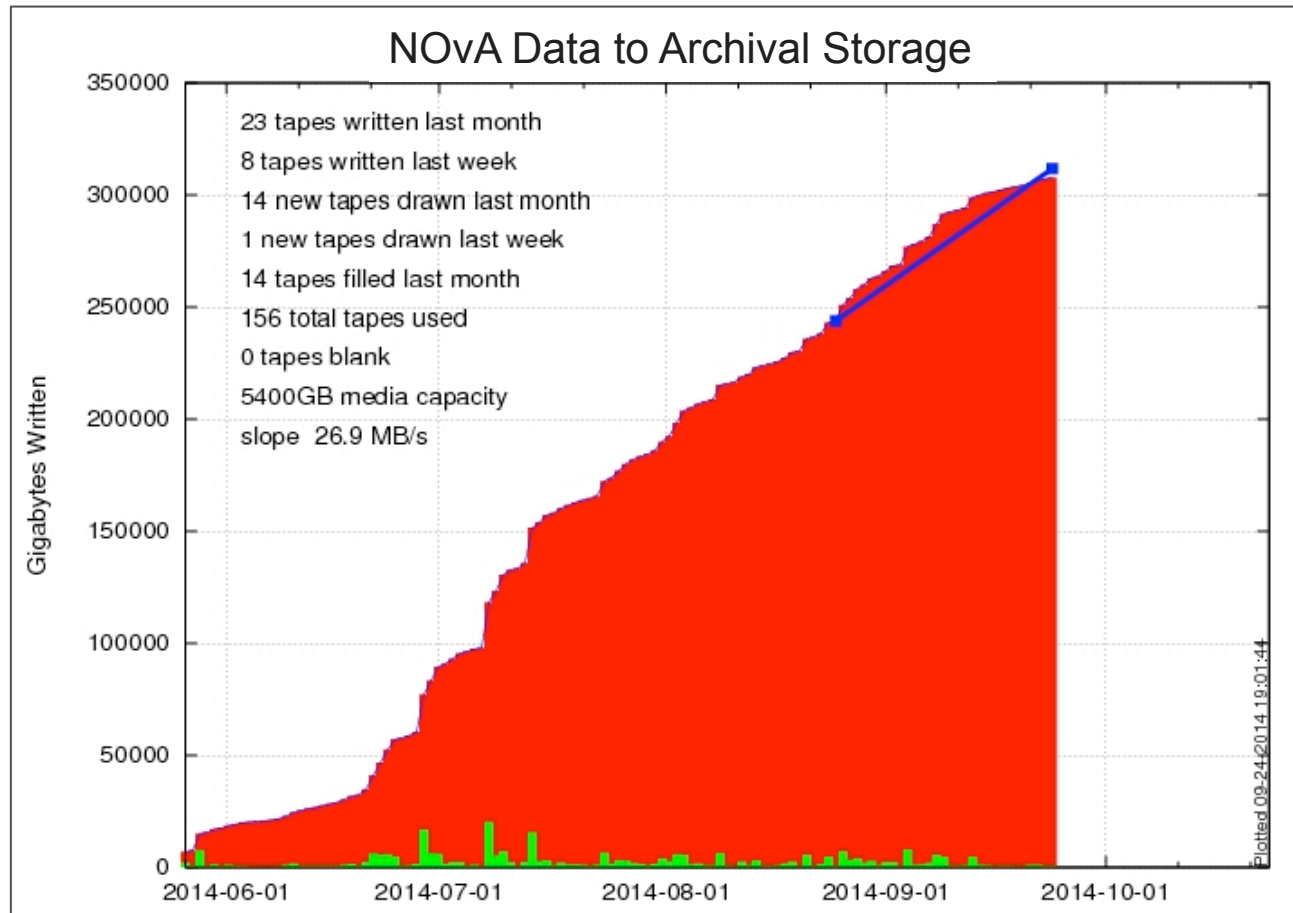
NO ν A (2014-)

Data Bandwidth

- How long does it take HEP to create 1 PB of data?

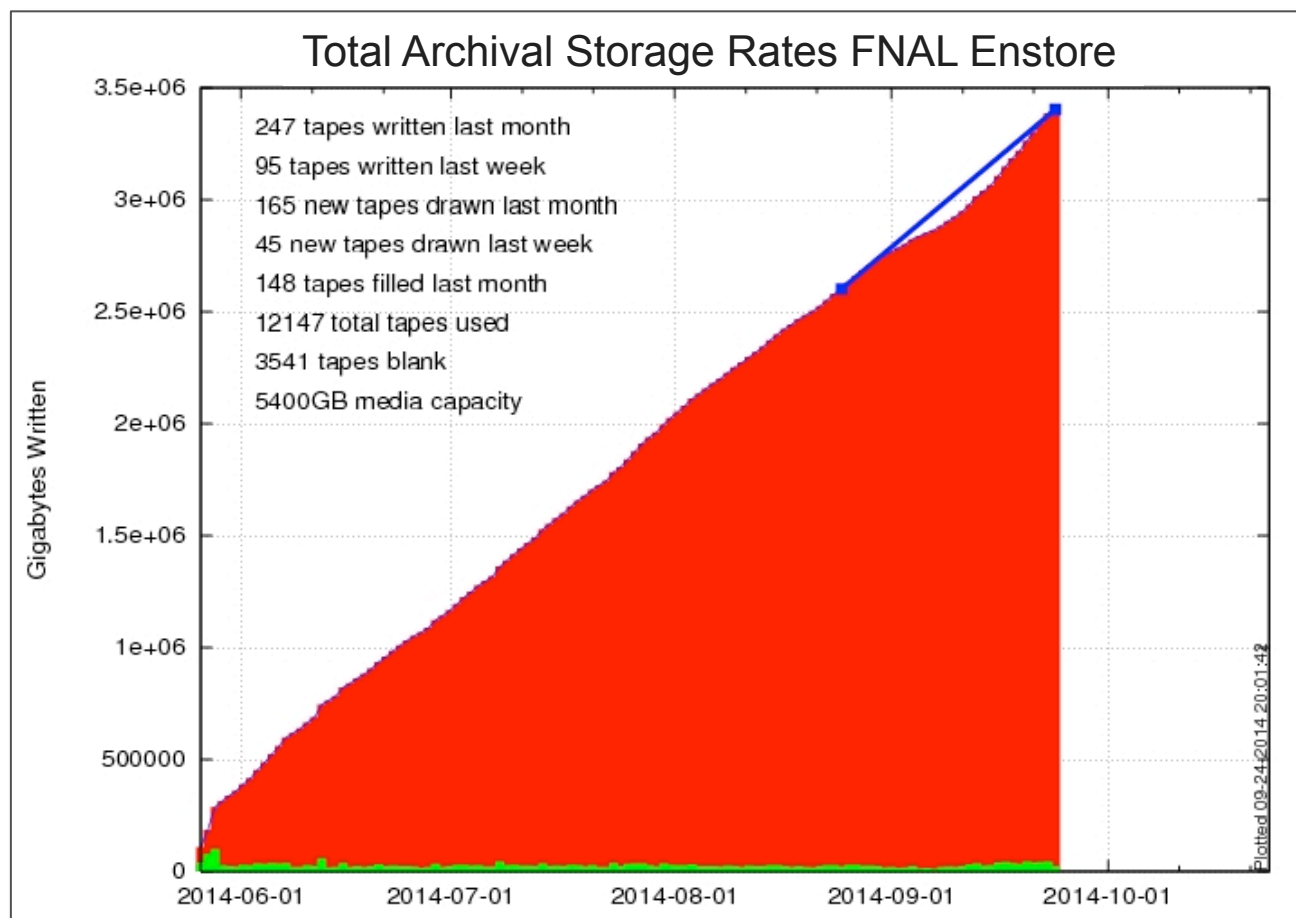
Epoch	Experiment	Type	Data Tier	Readout Time/PB
2001-11	DØ/CDF	Colliders, Fermilab	Final Data Set	8 PB/10 yr
2014-Present	Nova	Neutrino, Fermilab	Raw	90.8 s
			Into L3 farms	2.8 days
			to disk	214 days
2008-Present	CMS	Collider, LHC	Raw	12.5 s
			Into L3 farms	13.9 hr
			to disk	57.9 days

Data Volume – Single Experiment



- Data into FNAL Archival Storage for NOvA experiment ~125 TB/month
- Needs to be fully catalogged and retrievable to be of use

Data Volume – Aggregated Labwide



- **Total archival data volume (2014-09-25) 83509.685 TB**
- Data into archival storage lab wide ~1.3 PB/month
- Handled through mixed catalog SAM + CMS Data Management

Data Bandwidth

- How long does it take to transfer 1 TB/PB of data?

Network Speed	TB Transfer	PB Transfer	Network Class
100 GigE	88 sec	1 day	Research Backbones, Data center interconnects
10 GigE	14.8 min	1.5 weeks	Data center LANs
GigE	2.4 hr	103 days	HS WANs
100 Mbit	1.0 days	2.8 years	Commerical WANs
T1	102 days	183 years	Commercial WANs
300 Baud	117 years	117 centuries	my first modem

- Experiment Acquisition rates are limited by the ability to consume the data

The Problem

- There needs to be an efficient method of classifying, storing and retrieving vast amounts of sensor data
- Needs to be:
 - Capable of handling ANY type of data
 - Compatible with distributed analysis and analytics gathering computing models
 - Compatible with arbitrary data stores
 - Scalable (billions of rows)
 - Optimize access to the data

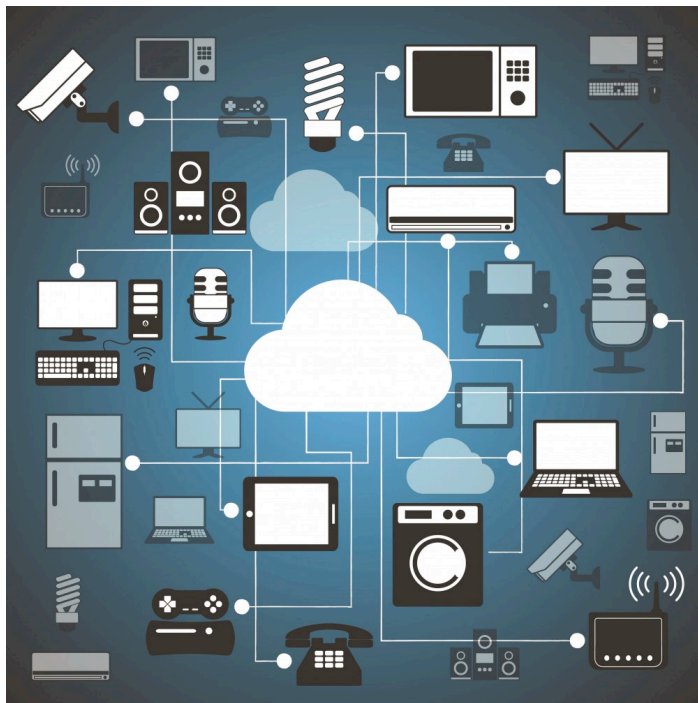
These the are the same problems faced by IoE

The Solution

- **Use the tools developed for High Energy Physics that are already tailored to perform large scale data analysis using object data stores indexed with metadata**
- **Perform analysis and mining operations on distributed clusters and computing clouds in the same manner that HEP analysis are done now**
- **Use the comprehensive data handling suite (SAM) along with extensions for Big Data**

IoE and HEP DAQ

- How is the “Internet of Things” like a neutrino experiment?

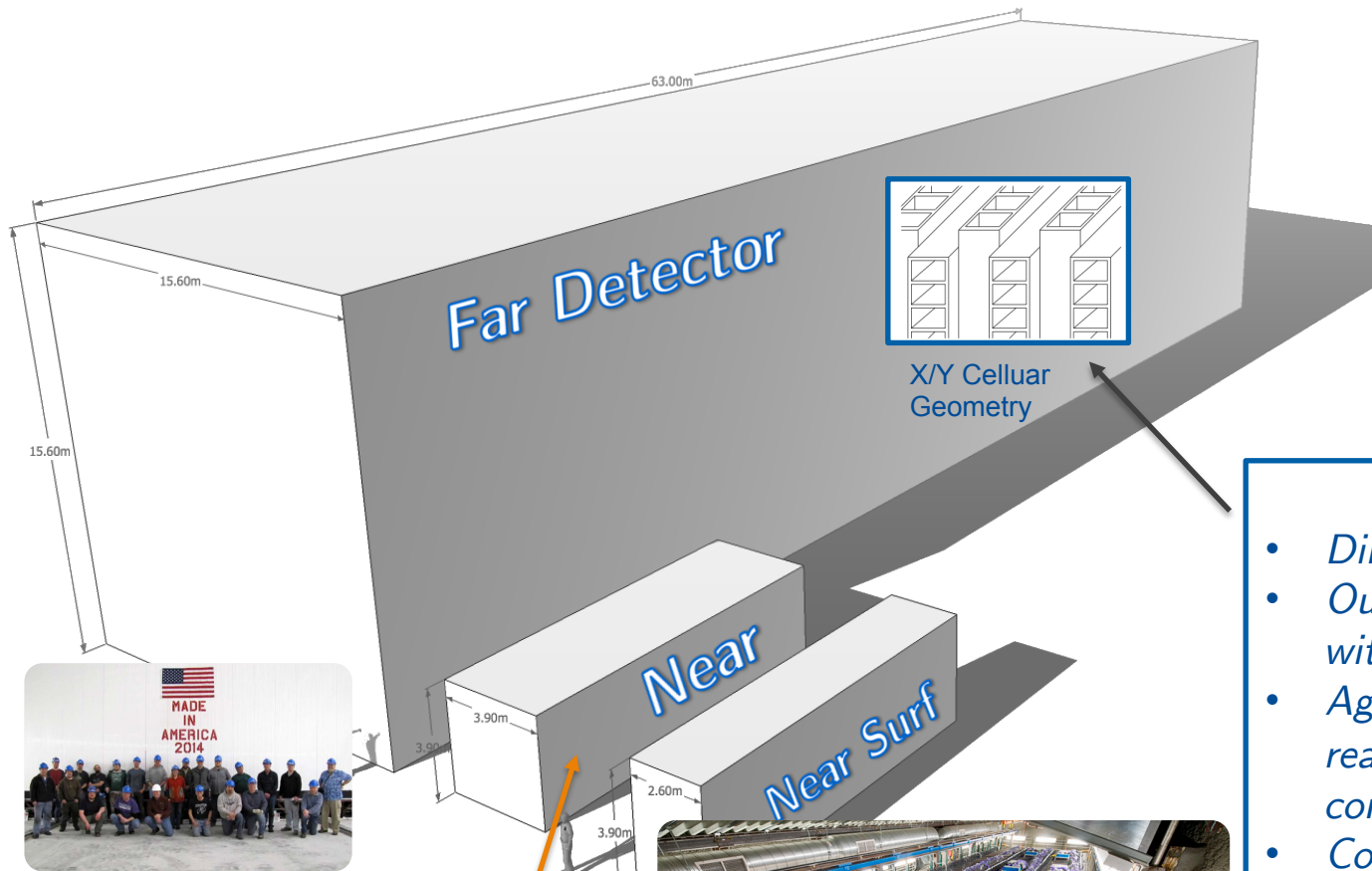


IoE and HEP DAQ

- How is the “Internet of Things” like a neutrino experiment?
- Current High Energy Physics experiments already use
 - high channel count (and data bandwidth) sensor packs
 - connect them through high speed networking
 - aggregate the data flows through layers of smart switches
 - analyze the data flows in real time to create “trigger” events
 - record vast data volumes in modern storage arrays
 - mine the recorded data in distributed, global computing resources

This is exactly what the IoE looks to do

NOvA Detectors



Far Detector

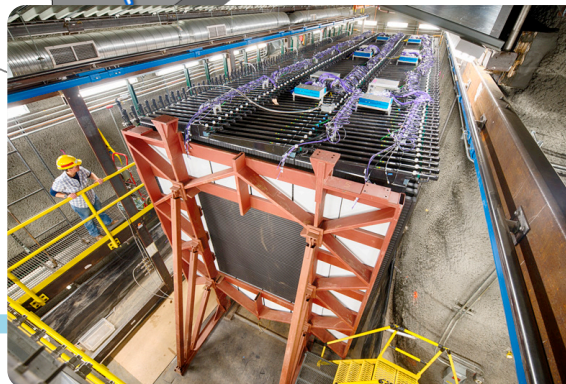
- Dims: 53x53x180 ft
- Outfitted on detector surfaces with over 344,000 sensors
- Aggregated through 10,752 readouts and 168 custom data combiners & smart switches
- Continuous readout system at 2M full detector frames/s
- Data bandwidth is 4.7 Gb/s
- Analyzed in realtime by compute farm to do triggering
- Triggered streams ≈ 1 TB/day
- $\approx 400:1$ data reduction

Nova Detectors

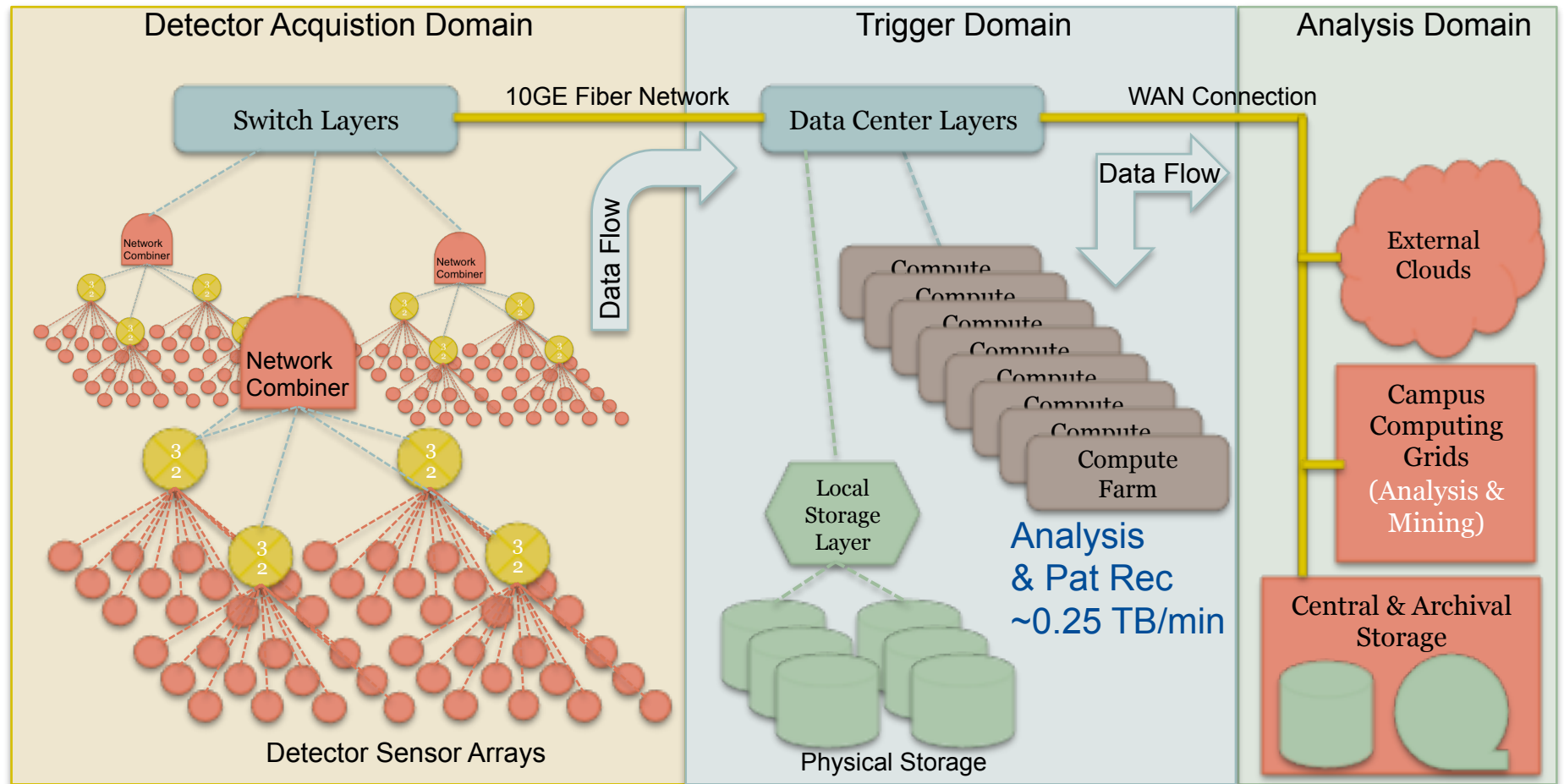
Flagships of the US neutrino program. Entered full operations Sept. 2014.

Accumulated and processed ≈ 0.75 PB of data during first six months of running

Big Data



HEP Data Acquisition

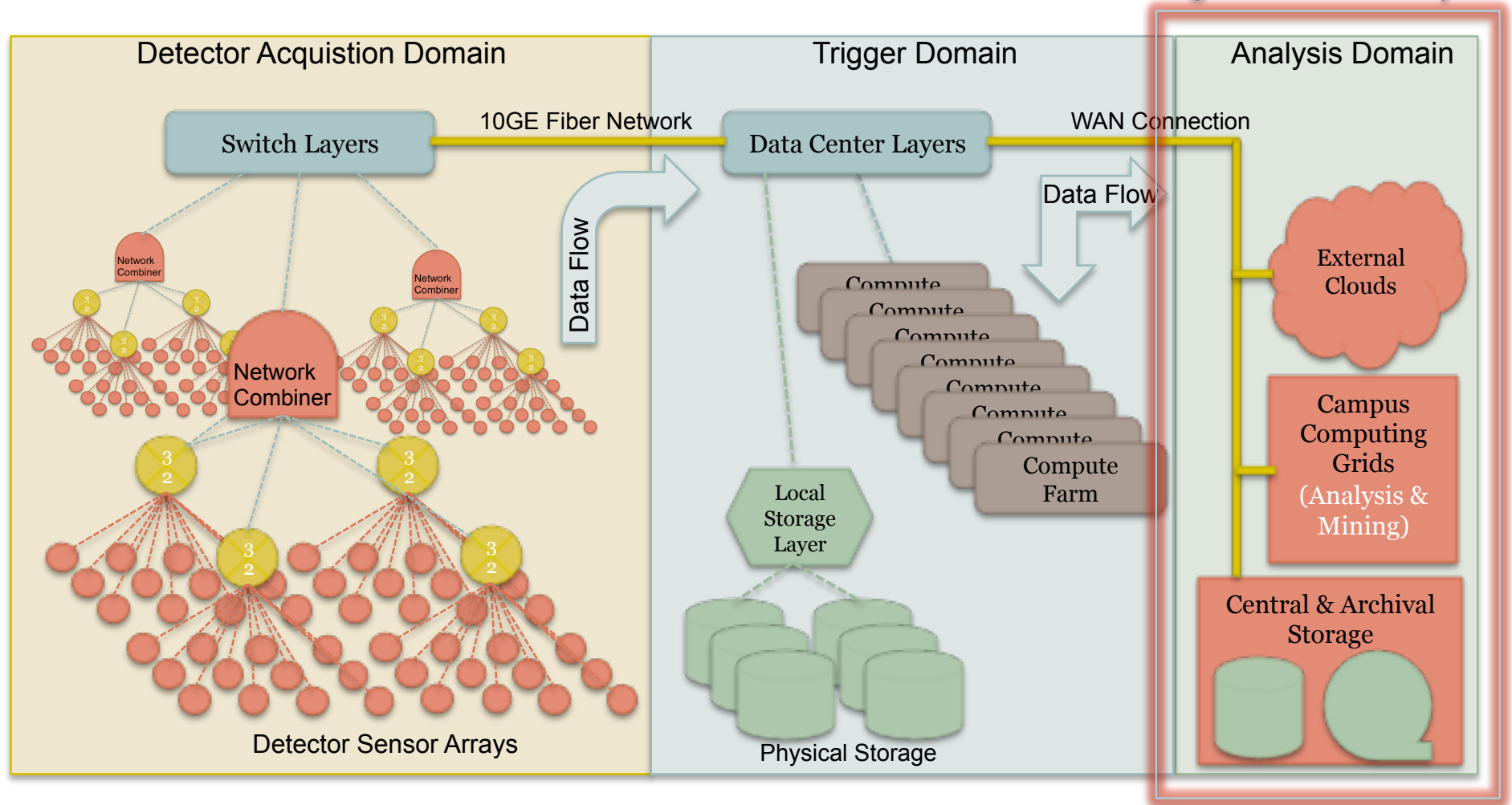


~500k sensor array cells
at 2-16 Mhz readout

Realtime analysis/trigger/extraction
Decision latencies (5ms-20min)

HEP Data Acquisition

Data Management & Analysis



Data Handling Components

- **S.A.M.**
 - Comprehensive Data Management
 - Metadata based catalog system
 - Replica Catalog
 - Project framework for organizing dataflow
- **Fermi File Transfer Service (F-FTS)**
 - Data registration and replication tool
 - Handles bulk transfers/cloning between sites
 - Module interface for data agnostic interface to SAM catalogs
- **Data Protocol Abstraction Layer (Data-PAL)**
 - Data transport abstraction layer
 - Handles protocol dependent authentication
- **Cern Virtual Machine File System (CVMFS)**
 - Distributed filesystem for delivery of application code

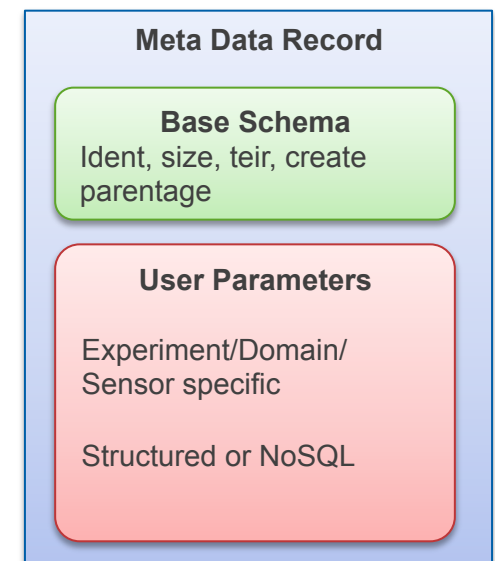
Sequential Access w/ Metadata (SAM)

- SAM is a model and application suite used for storing and mining large data sets for High Energy Physics.
 - Originates w/ Fermilab collider experiments
 - Expanded to manage data volumes of intensity frontier experiments
- Proven track record of cataloging and mining multi-petabyte datasets for the CDF and DØ experiments
(delivering/processing >1.5 PB/week for analysis)

SAM Statistics	
Actively Managed	230M files
Largest Catalog	180M files
Smallest Catalog	464K files
Peak delivery*	380 TB/day
Data Cataloged	>18 PB
Growth (projected 2014)	4 PB/yr

SAM Concept

- “Object based” data, replica and project catalog
- Each data object is registered in the catalog along with metadata describing it.
 - Two components to the metadata
 - **Base schema – General Object Information**
 - identifier, size, data tier, begin/end times, parentage/provenance
 - **User parameters – Data content specific fields**
 - Detector type, location, trigger stream, etc...
 - Only base schema is required
 - Simplifies registration of foreign/legacy data with catalog systems
- “**Datasets**” are then defined via queries against the meta data.
 - Evaluate to the set of objects to retrieve/analyze



SAM Concept (2)

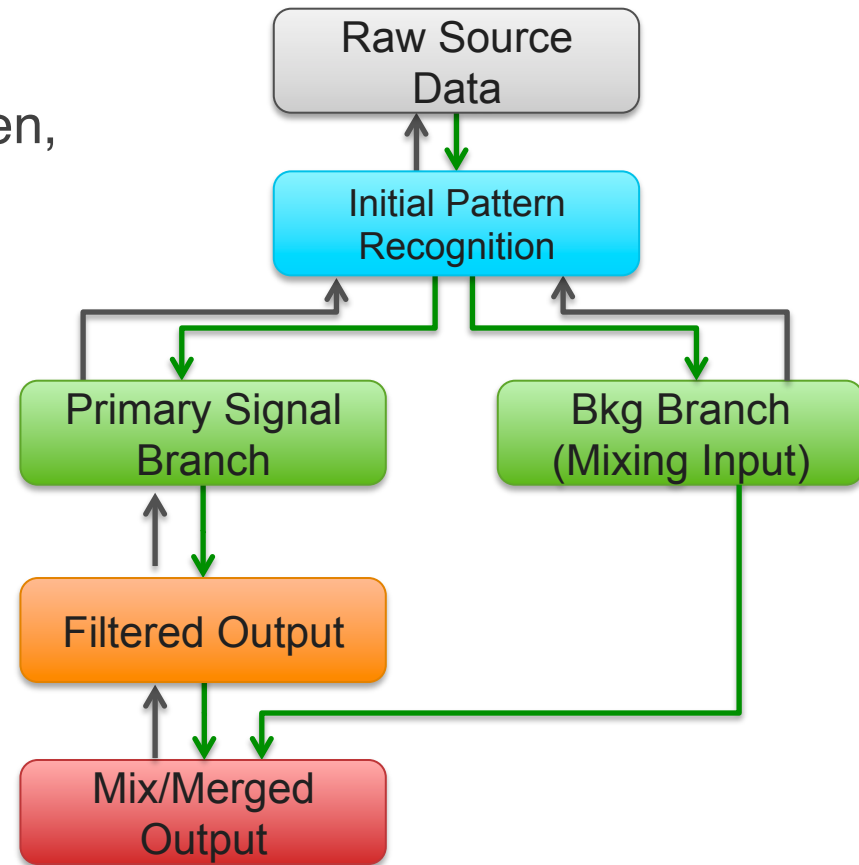
- SAM Analysis Projects provide a sophisticated management layer that optimizes scheduling and delivery of data from a dataset
 - Collections of compute resources register as clients of the project
 - Data objects from the dataset are assigned to the compute nodes
 - Ordering and dispatch of data is designed to take advantage of replica information [data locality, local cache states]
 - Retrieval from high latency storage (tape) are optimized and scheduled to maximize total throughput (i.e. caches are pre-populated with groups of data from sets)

SAM Concept (3)

- Analysis result objects are declared back to the SAM system along with new or expanded metadata including parentage information
- Creates full provenance chain that can be tracked or reproduced at any time.
- Allows for “self organizing” analysis chains
 - The metadata & parentage allows the processing chain to know what files still need to be run through a given algorithm
 - Allows for full automation of the analysis

Advanced Data Sets

- Use parentage to specify different types of complex relationships
 - Parents, children, grand children, peers, mixings, hybrids etc...
- Preserves the full parentage of each object
 - Fully traceable
 - Objects inherit meta-info
 - Permits containers
- Allows complex self organization and Identification
- Allows cross correlations in signal trees



SAM Dataset Language

- Datasets are defined using a “dimensions” language syntax against the metadata fields to construct a data query.

```
“(data_tier = raw) AND (start_time > 2014-09-25T12:00:00) AND (events > 0)  
AND ( (physics.signal_prob > 0.9) OR (physics.bkg_prob < 0.3) )”
```

- Designed to simplify the query and hide any underlying SQL schema details
 - Permits use of NoSQL like extensions to metadata
 - Query translation & optimization is performed server side
 - Guards and filters against ill-formed/defined requests
 - Scalable through frontend/backend mappings
 - Permits webAPI to be exploited for human interfaces to data stores

Tailored Web Interfaces

- Web Interfaces are dynamically tailored to the applications' data catalogs
 - Specific metadata and hierarchical relationships are exposed
 - Dimensions queries captured
- Provides the first stage in providing human interface

Catalog dynamically helps the user find the data

NOvA Monte Carlo Dataset Definition Editor

This page is designed to allow you to define your own custom data sets based on the current NOvA Monte Carlo data files that have been generated. To access the raw data or processed data set pages follow these links:
[Raw Data Files and Sets](#)
[Processed \(Reco\) Data Files and Sets](#)
For more information on creating and using custom data sets see:
[SAM Data Sets Wiki](#)

Monte Carlo Selection Criteria

Previously Defined Data Sets:
 Group/User: nova
 (To start with a previously defined dataset)

= to
 =
 =
 =
 =
 =
 =
 =
 >
 = (Example: 'Geometry/gdml/ndos.gdml')
 = (Example: cospics_ndos_10000_r1_99.fcl)
 =
 (Date format: 2011-05-09 or Date/Time format: 2011-05-09T23:46:04)

Logical Operators
Use these operators to join your criteria together.

Data Set Definition (Dimensions query):
(you may also edit this query string directly to add custom fields to your query)

 (SAM Translate)

Name your dataset: Save As: user: group:
Datasets can have an arbitrary name but should not include spaces or special characters (underscores and dashes are permitted)

SAMWeb

- Modern http based Client/Server tools
- Simplifies client access to SAM functionality
 - Eliminates the need for dedicated SAM stations at sites
 - Allows experiments universal access to SAM resources from non-FNAL locations
 - Allows cross platform access to the SAM toolset (Linux/Unix, OSX, anything that can run Python or talk http)
- Improves upon the functions/tasks people really use
 - Simplified function calls
 - Optimizations to common tasks (i.e. multi-file and bulk operations)

SAM Replica Catalog

- Replica catalog contains a 1-to-many mapping of the data element's identity to the different storage locations where it resides
- Supports
 - Durable locations – central storage elements, cache disks
 - Archival locations – tape systems and other high latency storage
 - Virtual locations – declarations of file locality that don't yet have physical representation (i.e. files in transit between sites)
- Namespaces are kept distinct between SAM catalogs
- Mappings preserve the namespaces of the actual storage resources

SAM Replica Catalog

- Data elements are mapped logically to the storage resource
- Access protocols supported by a storage resource are registered with the catalog system
- Provides a full mapping of:
 data element → storage element → access URI
- Allows for seamless multi-protocol data access

Replica Example

- Simple locate of a data element:

```
# Locate a data element  
samweb locate-file <data element ID>
```

Multiple replicas exist:

- central disk
- tape + dCache

Filename serves as identifier

```
samweb locate-file ndos_r0016356_s00_ddtTriCell.raw
```

Returns:

```
novadata:/nova/data/rawdata/NDOS/000163/16356/ddttricell
```

```
enstore:/pnfs/nova/rawdata/NDOS/runs/000163/16356(3618@vpl001)
```

resource

namespace path

volume/tape label

- Access is then requested via some protocol:

Filename serves as identifier

```
samweb get-file-access-url -schema=xroot ndos_r0016356_s00_ddtTriCell.raw
```

```
xroot://ndca1.fnal.gov:1094/pnfs/fnal.gov/user/nova/rawdata/NDOS/runs/000163/16356/ndos_r00016356_s00_ddtTriCell.raw
```

streaming supported only
from dCache location

```
samweb get-file-access-url -schema=gsiftp ndos_r0016356_s00_ddtTriCell.raw
```

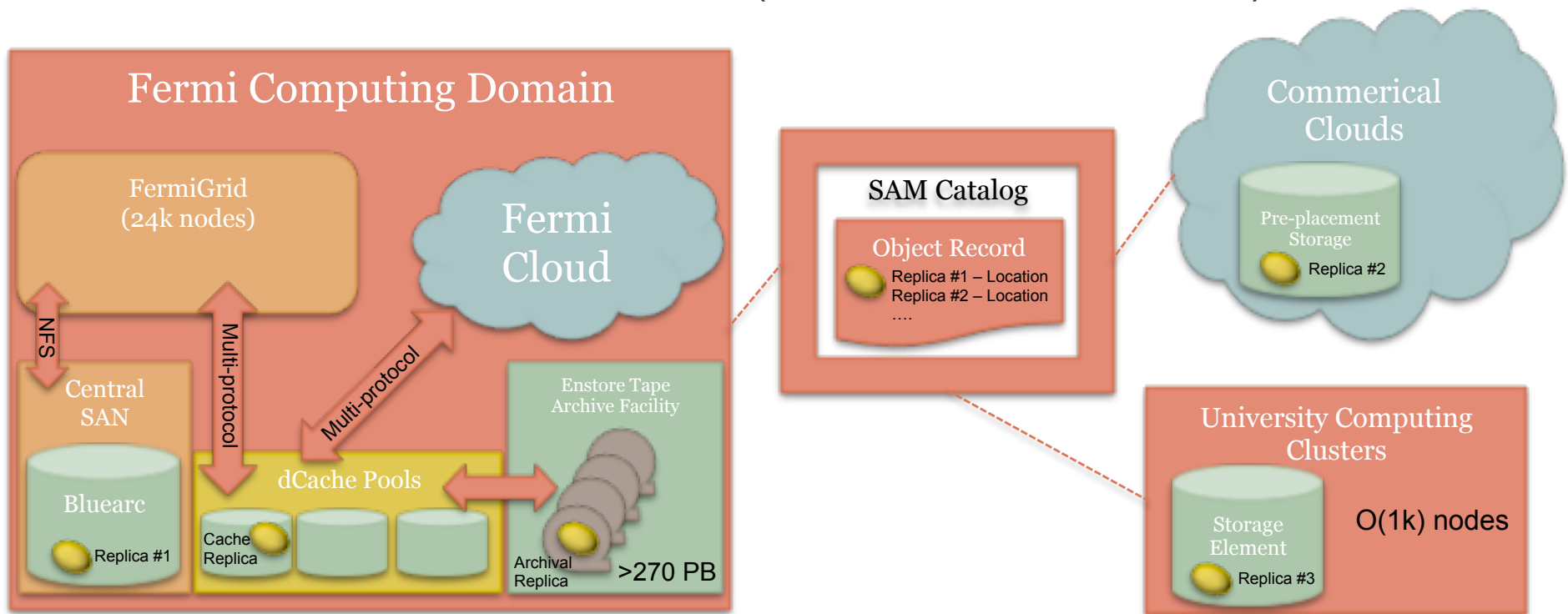
```
gsiftp://fg-bestman1.fnal.gov:2811/nova/data/rawdata/NDOS/000163/16356/ddttricell/ndos_r00016356_s00_ddtTriCell.raw  
gsiftp://ndca1.fnal.gov:2811/rawdata/NDOS/runs/000163/16356/ndos_r00016356_s00_ddtTriCell.raw
```

gsiftp support from both
locations

All supported locations for a protocol are returned

Replicas with locality

- Data elements are placed on storage “close” to a given compute cluster
 - SAM returns close location (with fall back to others)



Allows for integration of multiple storage technologies across sites

Data Staging

- SAM is designed to optimize access to high latency storage through staging requests
 - Pre-staging mechanism to preemptively load large datasets
 - Primarily used to pre-populate a storage system with a large data set
 - Standard staging mechanism performs predictive load requests during project execution (i.e. requests for the next N data elements are dispatched in advance of the need for the elements)
 - Can utilize various types of cache systems
 - SAM managed cache disk
 - Externally managed cache systems (e.g. dCache)
 - Flexibility regarding how to deal with high latency systems

Data Replication & Cloning

- Specialized tool for cloning large data sets across WAN:

Fermi File Transfer Service (F-FTS)

- Designed for large scale replication of critical (irreplaceable) data
- Handles replication and organization of objects in the destination's namespace
 - Rule based replication (supports multiple replications per object)
- I/O & Bandwidth management
- Validates all transfers and performs data integrity checks
- Registers new locations with catalog
- Proven record of replicating over 6.3M files between Minnesota woods and Fermilab over commercial WAN

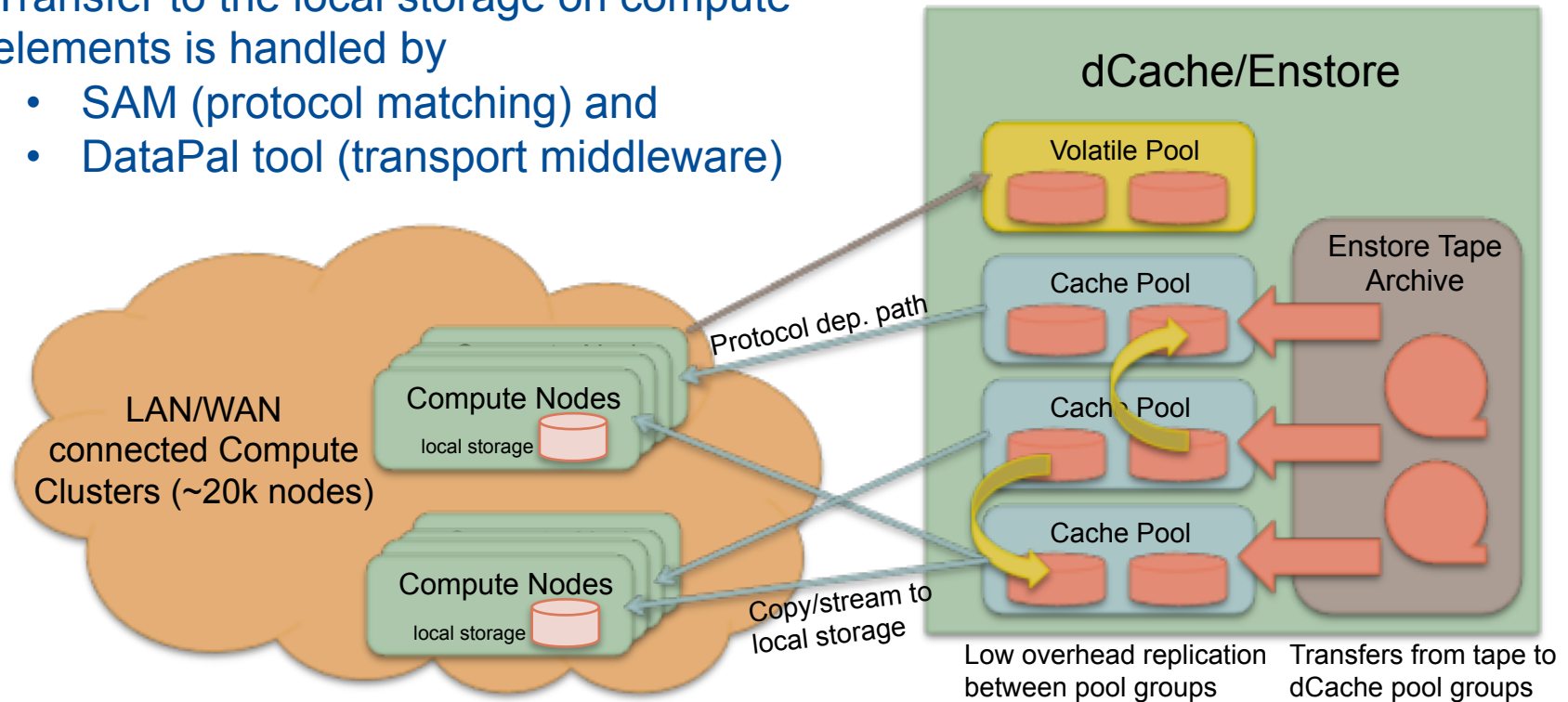


SAM Cache Characteristics

- Optimize access to tape: Groups and schedules tape access to minimize mounts
- Protects tape access from overload
 - Throttles access by cache node and project
- Protects cache system from overload
 - Throttles access to cache disks by node & project
- Protects cache from short time to live (TTL)
 - Prevents high throughput/turn over production activities from ejecting specialized skims for analysis
- File placement and expiration is automatic with LRU cache
 - Never pre-place/pin data in cache
 - Regular usage brings files in from tape
 - Popular data remains in cache
 - Old or unpopular data expires
 - Requires no human effort
- Automatic prefetching: SAM populates cache prior to job start and while jobs are busy
- Supports many transfer/storage protocols:
 - Utilize many different types of storage elements
 - Allows for opportunistic running at remote sites

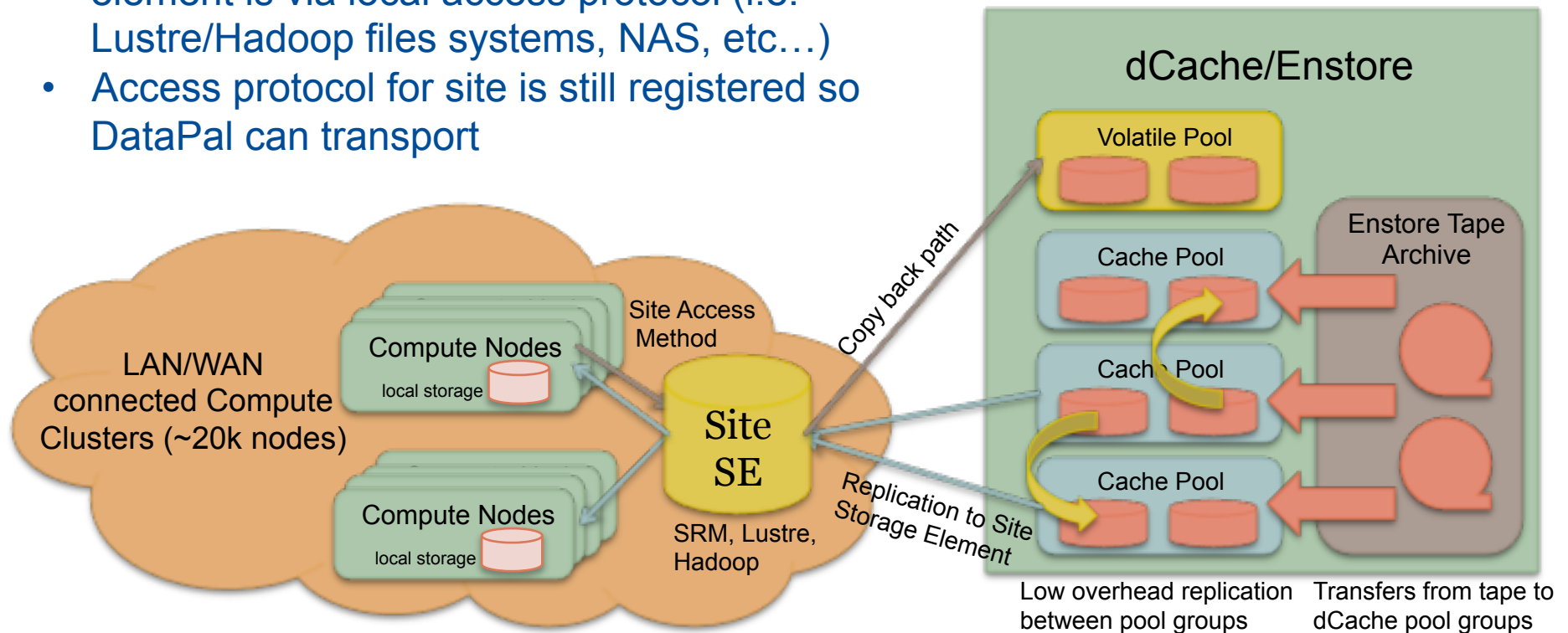
SAM and dCache Integration

- SAM acts as a scheduler, optimizer and throttle for the requested data streams
- dCache manages the pools groups and internal replications
- Transfer to the local storage on compute elements is handled by
 - SAM (protocol matching) and
 - DataPal tool (transport middleware)



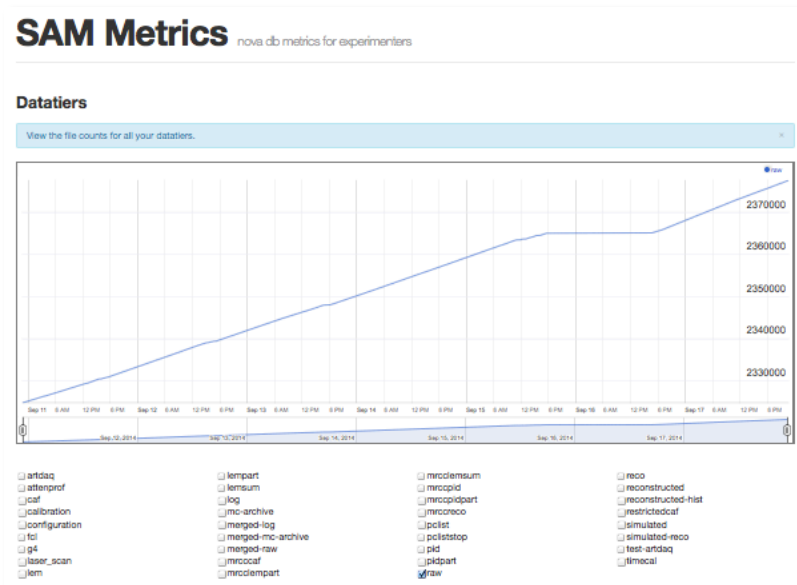
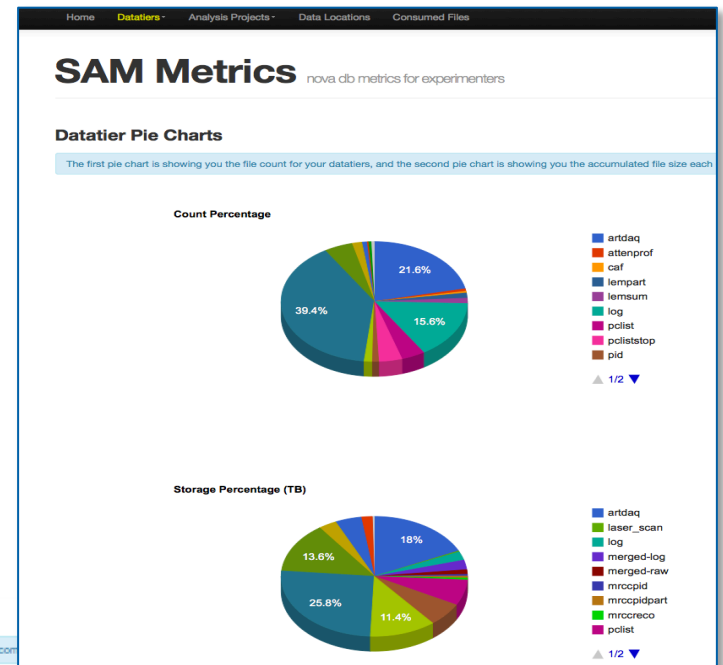
SAM and dCache Integration

- For remote sites we can insert intermediate cache stations close to the computing
- Replication is first to the local SE
- Subsequent streaming/copy to the compute element is via local access protocol (i.e. Lustre/Hadoop files systems, NAS, etc...)
- Access protocol for site is still registered so DataPal can transport



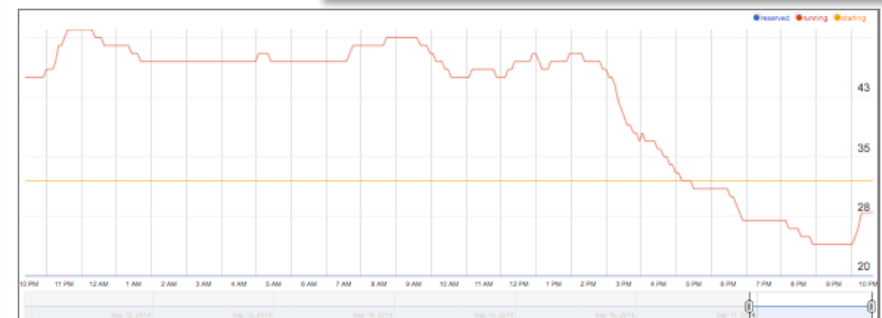
System Analytics

- Full (customizable) analytics and monitoring for each SAM instance
 - Access to data counts, storage footprints, recent activity
 - Tracks evolution of the system as data is accumulated



Analysis Projects

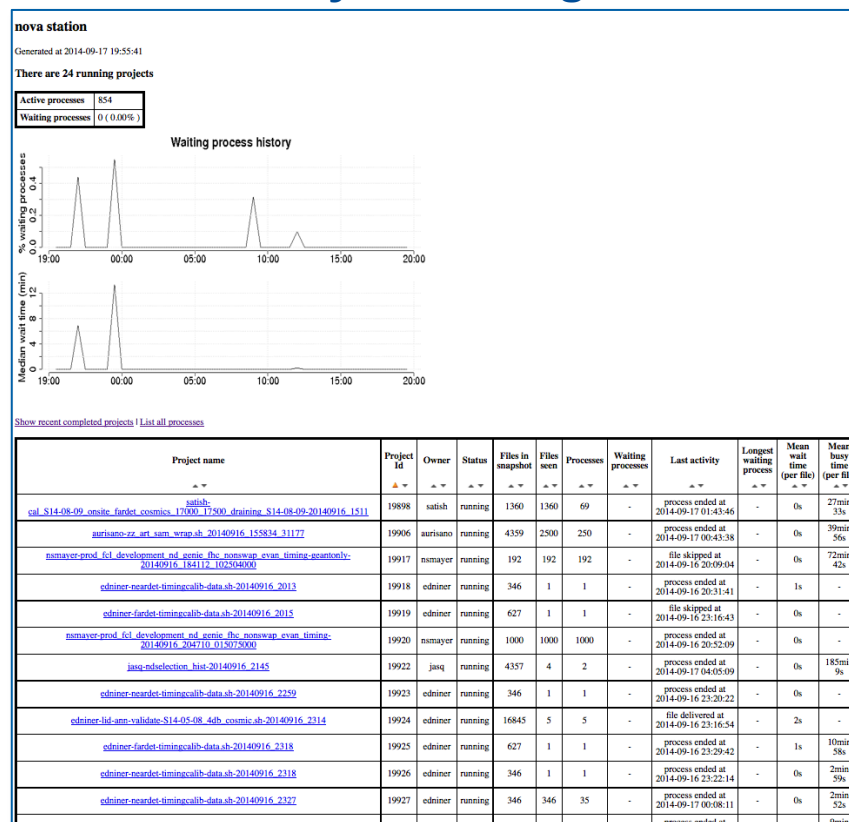
This chart shows you the accumulated total of complete/incomplete analysis projects.



System Monitoring

- Monitoring and status interfaces provide detailed drill down on analysis activities

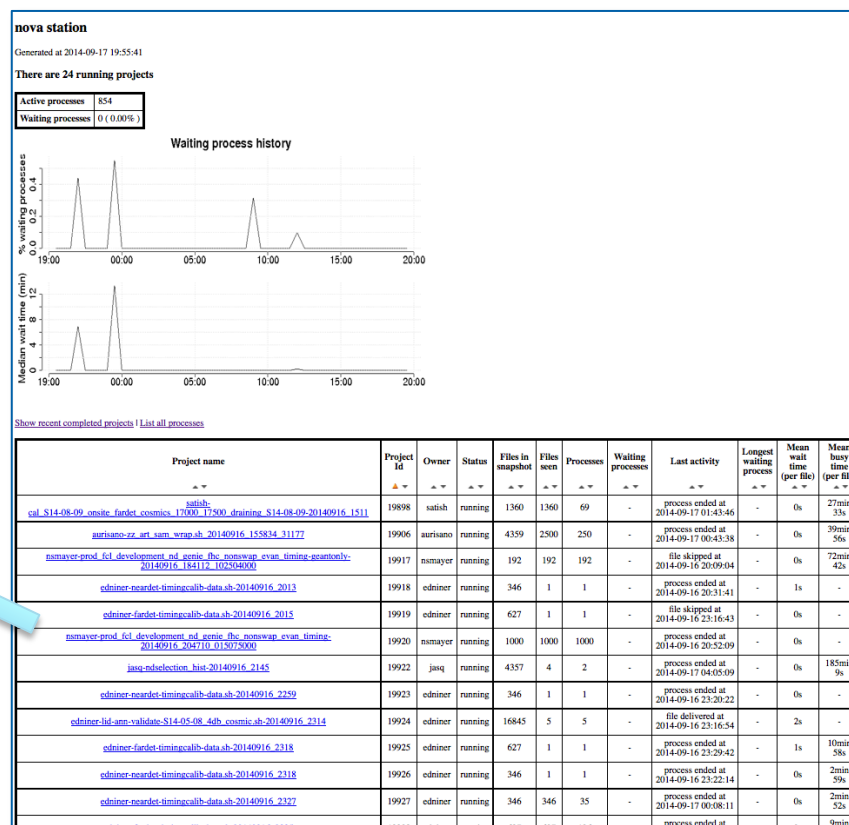
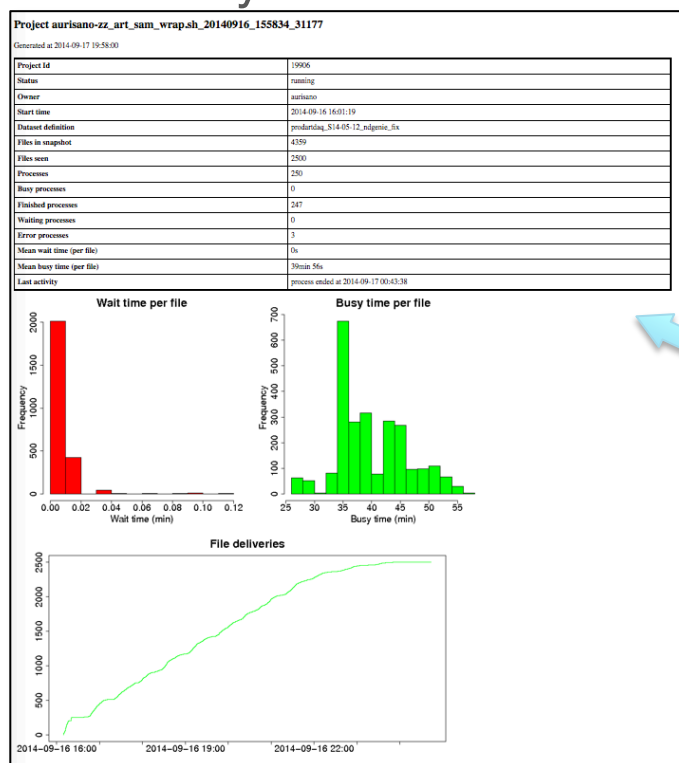
SAM Analysis Management



Monitoring of Global Project Activity

System Monitoring

- Monitoring and status interfaces provide detailed drill down on analysis activities

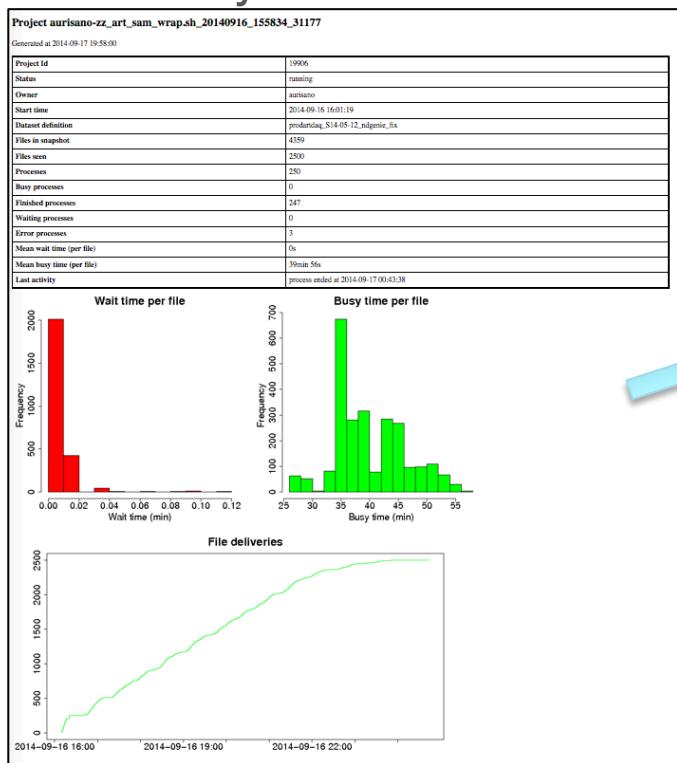


Monitoring of Global Project Activity

Provides detailed audit on data access, delivery and consumption

System Monitoring

- Monitoring and status interfaces provide detailed drill down on analysis activities



SAM + Grid Integration

Process id	Node name	Status	Description	Files seen	Last change	Waiting for	Mean wait time (per file)	Mean busy time (per file)
1398962	lfp0035.fnl.gov	completed	17668580.0	10	2014-09-16 23:24:48 (process ended - completed)	-	0s	43min 31s
1398963	lfp0204.fnl.gov	completed	17668572.0	10	2014-09-16 20:50:13 (process ended - completed)	-	0s	28min 4s
1398964	lfp0204.fnl.gov	completed	17668574.0	10	2014-09-16 20:50:13 (process ended - completed)	-	0s	28min 7s
1398965	lfp0006.fnl.gov	completed	17668575.0	10	2014-09-16 22:06:48 (process ended - completed)	-	0s	35min 43s
1398966	lfp0003.fnl.gov	completed	17668573.0	10	2014-09-16 22:19:38 (process ended - completed)	-	0s	37min 0s
1398967	lfp03072.fnl.gov	completed	17668582.0	10	2014-09-16 22:03:48 (process ended - completed)	-	0s	35min 25s
1398968	lfp0309.fnl.gov	completed	17668587.0	10	2014-09-16 20:45:41 (process ended - completed)	-	0s	27min 56s
1398969	lfp0025.fnl.gov	completed	17668590.0	10	2014-09-16 23:26:34 (process ended - completed)	-	0s	43min 41s
1398970	lfp0216.fnl.gov	completed	17668577.0	10	2014-09-16 20:57:10 (process ended - completed)	-	0s	28min 45s
1398971	lfp0018.fnl.gov	completed	17668588.0	10	2014-09-16 22:05:48 (process ended - completed)	-	0s	35min 37s
1398972	lfp0019.fnl.gov	completed	17668583.0	10	2014-09-16 23:27:18 (process ended - completed)	-	0s	44min 06s
1398973	lfp0044.fnl.gov	completed	17668594.0	10	2014-09-16 23:27:22 (process ended - completed)	-	0s	43min 06s
1398974	lfp0508.fnl.gov	completed	17668596.0	10	2014-09-16 22:33:09 (process ended - completed)	-	0s	38min 21s
1398975	lfp0013.fnl.gov	completed	17668591.0	10	2014-09-16 22:06:09 (process ended - completed)	-	0s	35min 39s
1398976	lfp0500.fnl.gov	completed	17668595.0	10	2014-09-16 22:44:27 (process ended - completed)	-	0s	39min 29s
1398977	lfp0050.fnl.gov	completed	17668595.0	10	2014-09-16 23:31:00 (process ended - completed)	-	0s	44min 8s
1398978	lfp0059.fnl.gov	completed	17668599.0	10	2014-09-16 23:27:25 (process ended - completed)	-	0s	43min 46s
1398979	lfp0010.fnl.gov	completed	17668581.0	10	2014-09-16 23:01:17 (process ended - completed)	-	0s	35min 9s
1398980	lfp004034.fnl.gov	completed	17668604.0	10	2014-09-17 00:28:55 (process ended - completed)	-	0s	49min 35s
1398981	lfp0747.fnl.gov	completed	17668609.0	10	2014-09-16 23:04:33 (process ended - completed)	-	0s	35min 38s
1398982								

Process 1398964 [17668574.0]
Generated at 2014-09-17 20:08:00

Process id	1398964
Node name	lfp0204.fnl.gov
Status	completed
Start time	2014-09-16 16:09:29
Files seen	10
Last activity	2014-09-16 20:50:48 (process ended - completed)

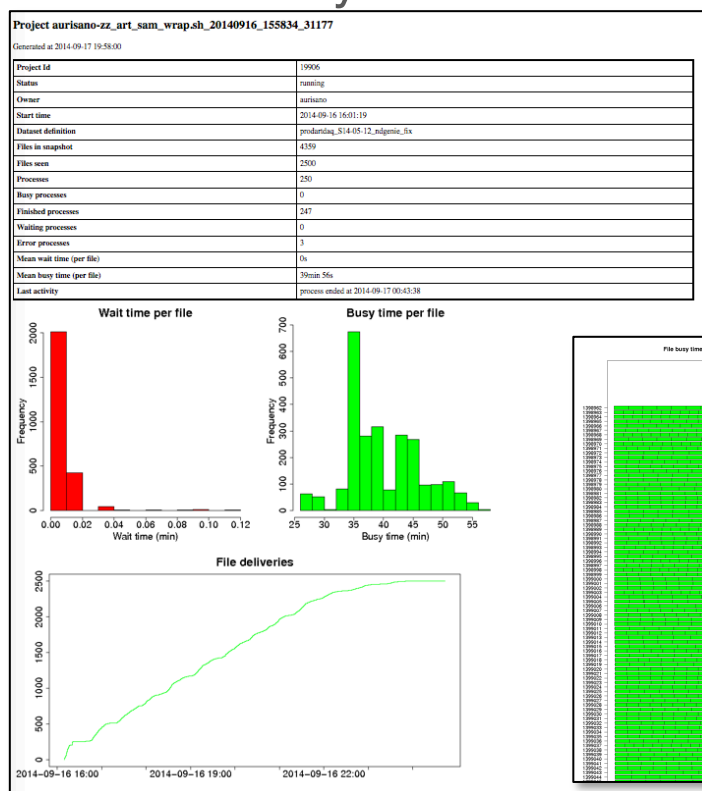
File id	File name	Status	Delivered	Closed	Waited for	Busy for
0015780	nsarlet_gmtc_fbc_nonswap_1000_010000041_s57_S14-05-12_v3_20140819_003614.fnl.gov_140845360_4775_0.sim.overlay.dag.root	consumed	2014-09-16 16:09:29	2014-09-16 16:38:09	0s	28min 40s
0015460	nsarlet_gmtc_fbc_nonswap_1000_010000054_s25_S14-05-12_v3_20140819_003614.fnl.gov_140845048_13443_0.sim.overlay.dag.root	consumed	2014-09-16 16:38:09	2014-09-16 17:05:26	0s	27min 17s
0015740	nsarlet_gmtc_fbc_nonswap_1000_010000041_s43_S14-05-12_v3_20140819_003614.fnl.gov_140845876_48792_0.sim.overlay.dag.root	consumed	2014-09-16 17:05:26	2014-09-16 17:34:35	0s	28min 9s
0015482	nsarlet_gmtc_fbc_nonswap_1000_010000044_s23_S14-05-12_v3_20140819_003614.fnl.gov_1408461784_38345_0.sim.overlay.dag.root	consumed	2014-09-16 17:34:35	2014-09-16 18:02:43	0s	28min 8s
0015480	nsarlet_gmtc_fbc_nonswap_1000_010000046_s21_S14-05-12_v3_20140819_003614.fnl.gov_140846760_10405_0.sim.overlay.dag.root	consumed	2014-09-16 18:02:43	2014-09-16 18:30:21	1s	27min 37s
0015596	nsarlet_gmtc_fbc_nonswap_1000_010000038_s41_S14-05-12_v3_20140819_003614.fnl.gov_1408461483_53098_0.sim.overlay.dag.root	consumed	2014-09-16 18:30:21	2014-09-16 18:57:53	0s	27min 12s
0015116	nsarlet_gmtc_fbc_nonswap_1000_010000090_s06_S14-05-12_v3_20140819_003614.fnl.gov_1408467225_31931_0.sim.overlay.dag.root	consumed	2014-09-16 18:57:53	2014-09-16 19:26:22	0s	28min 49s
0015972	nsarlet_gmtc_fbc_nonswap_1000_010000072_s47_S14-05-12_v3_20140819_003614.fnl.gov_140846779_16775_0.sim.overlay.dag.root	consumed	2014-09-16 19:26:22	2014-09-16 19:53:47	0s	27min 25s
0015127	nsarlet_gmtc_fbc_nonswap_1000_010000011_s20_S14-05-12_v3_20140819_003614.fnl.gov_1408469306_47566_0.sim.overlay.dag.root	consumed	2014-09-16 19:53:47	2014-09-16 20:22:03	0s	28min 16s
0017722	nsarlet_gmtc_fbc_nonswap_1000_010000078_s18_S14-05-12_v3_20140819_003614.fnl.gov_1408454587_31927_0.sim.overlay.dag.root	consumed	2014-09-16 20:22:03	2014-09-16 20:50:48	1s	28min 44s

Tracks all data down to the compute nodes and process level

Provides detailed audit on data access, delivery and consumption

System Monitoring

- Monitoring and status interfaces provide detailed drill down on analysis activities



Process id	Node name	Status	Description	Files seen	Last change	Waiting for	Mean wait time (per file)	Mean busy time (per file)
1398962	lfp0035.fnl.gov	completed	17868580.0	10	2014-09-16 23:24:48 (process ended - completed)	-	0s	43min 31s
1398963	lfp0204.fnl.gov	completed	17868572.0	10	2014-09-16 20:50:13 (process ended - completed)	-	0s	28min 4s
1398964	lfp0204.fnl.gov	completed	17868574.0	10	2014-09-16 20:50:48 (process ended - completed)	-	0s	28min 7s
1398965	lfp0006.fnl.gov	completed	17868575.0	10	2014-09-16 22:06:48 (process ended - completed)	-	0s	35min 43s
1398966	lfp0003.fnl.gov	completed	17868573.0	10	2014-09-16 22:19:38 (process ended - completed)	-	0s	37min 0s
1398967	lfp0022.fnl.gov	completed	17868582.0	10	2014-09-16 22:03:48 (process ended - completed)	-	0s	35min 25s
1398968	lfp0209.fnl.gov	completed	17868587.0	10	2014-09-16 20:45:01 (process ended - completed)	-	0s	27min 0s
1398969	lfp0025.fnl.gov	completed	17868590.0	10	2014-09-16 23:26:34 (process ended - completed)	-	0s	43min 41s
1398970	lfp0216.fnl.gov	completed	17868577.0	10	2014-09-16 20:57:10 (process ended - completed)	-	0s	28min 45s
1398971	lfp0018.fnl.gov	completed	17868588.0	10	2014-09-16 22:05:48 (process ended - completed)	-	0s	35min 37s
1398972	lfp0019.fnl.gov	completed	17868583.0	10	2014-09-16 23:37:18 (process ended - completed)	-	0s	44min 0s
1398973	lfp0044.fnl.gov	completed	17868594.0	10	2014-09-16 23:27:22 (process ended - completed)	-	0s	43min 0s
1398974	lfp0208.fnl.gov	completed	17868596.0	10	2014-09-16 22:33:09 (process ended - completed)	-	0s	38min 21s
1398975	lfp0013.fnl.gov	completed	17868599.0	10	2014-09-16 22:06:09 (process ended - completed)	-	0s	35min 39s
1398976	lfp0200.fnl.gov	completed	17868591.0	10	2014-09-16 22:44:27 (process ended - completed)	-	0s	39min 29s
1398977	lfp0050.fnl.gov	completed	17868595.0	10	2014-09-16 23:31:00 (process ended - completed)	-	0s	44min 8s
1398978	lfp0059.fnl.gov	completed	17868599.0	10	2014-09-16 23:27:25 (process ended - completed)	-	0s	43min 0s
1398979	lfp0010.fnl.gov	completed	17868581.0	10	2014-09-16 22:01:17 (process ended - completed)	-	0s	35min 9s
1398980	lfp004034.fnl.gov	completed	17868604.0	10	2014-09-17 00:28:55 (process ended - completed)	-	0s	49min 35s
1398981	lfp0047.fnl.gov	completed	17868605.0	10	2014-09-16 23:04:33 (process ended - completed)	-	0s	35min 38s
1398982								

Process 1398964 [17868574.0]
Generated at 2014-09-17 20:08:00

Process id	1398964
Node name	lfp0204.fnl.gov
Status	completed
Start time	2014-09-16 16:09:29
Files seen	10
Last activity	2014-09-16 20:50:48 (process ended - completed)

File id	File name	Status	Delivered	Closed	Waited for	Busy for
0015780	ncardet_gene_bsc_novseq_1000_010000041_07_S14-05-12_v3_20140819_003614.fnl.gov_1408453460_4775_0.sim.overlay.dag.root	consumed	2014-09-16 16:09:29	2014-09-16 16:38:09	0s	28min 40s
0015781	ncardet_gene_bsc_novseq_1000_010000054_025_S14-05-12_v3_20140819_003614.fnl.gov_1408455048_13443_0.sim.overlay.dag.root	consumed	2014-09-16 16:38:09	2014-09-16 17:05:26	0s	27min 17s
0015782	ncardet_gene_bsc_novseq_1000_010000041_043_S14-05-12_v3_20140819_003614.fnl.gov_1408458756_48792_0.sim.overlay.dag.root	consumed	2014-09-16 17:05:26	2014-09-16 17:34:35	0s	28min 9s
0015783	ncardet_gene_bsc_novseq_1000_010000044_021_S14-05-12_v3_20140819_003614.fnl.gov_1408461784_38343_0.sim.overlay.dag.root	consumed	2014-09-16 17:34:35	2014-09-16 18:02:02	1s	28min 3s
0015784	ncardet_gene_bsc_novseq_1000_010000046_021_S14-05-12_v3_20140819_003614.fnl.gov_1408461960_10405_0.sim.overlay.dag.root	consumed	2014-09-16 18:02:02	2014-09-16 18:30:21	1s	27min 37s
0015785	ncardet_gene_bsc_novseq_1000_010000038_041_S14-05-12_v3_20140819_003614.fnl.gov_1408461483_53098_0.sim.overlay.dag.root	consumed	2014-09-16 18:30:21	2014-09-16 18:57:53	0s	27min 12s
0015786	ncardet_gene_bsc_novseq_1000_010000090_026_S14-05-12_v3_20140819_003614.fnl.gov_1408467225_31931_0.sim.overlay.dag.root	consumed	2014-09-16 18:57:53	2014-09-16 19:26:22	0s	28min 49s
0015787	ncardet_gene_bsc_novseq_1000_010000072_047_S14-05-12_v3_20140819_003614.fnl.gov_1408467779_16775_0.sim.overlay.dag.root	consumed	2014-09-16 19:26:22	2014-09-16 19:53:47	0s	27min 25s
0015788	ncardet_gene_bsc_novseq_1000_010000011_020_S14-05-12_v3_20140819_003614.fnl.gov_1408469306_47566_0.sim.overlay.dag.root	consumed	2014-09-16 19:53:47	2014-09-16 20:22:03	0s	28min 16s
0015789	ncardet_gene_bsc_novseq_1000_010000078_018_S14-05-12_v3_20140819_003614.fnl.gov_1408454587_31927_0.sim.overlay.dag.root	consumed	2014-09-16 20:22:03	2014-09-16 20:50:48	1s	28min 44s

Tracks all data down to the compute nodes and process level

Provides allows for automated recovery of failed processes and files

Data Delivery – Data Protocol Abstraction Layer (Data-PAL)

- Swiss army knife of file delivery
- Designed to be a lightweight toolkit to handle the last leg of file delivery (to compute node)
 - “Smart” broker with location awareness
 - Integrated with SAM data catalogs
 - Modular system for transfer protocols
 - Supports gridftp, SRM, custom access protocols seamlessly
 - Provides single end user interface and syntax
 - Allows for workflows with “mixed” transport requirements
 - Handles authentication and certificate generation for FNAL users
 - Bidirectional operation (i.e. copy-in and copy-out)
 - Includes bulk copy operations
- Full abstraction layer
 - Most end users only need Data-PAL and aren’t exposed to protocol details

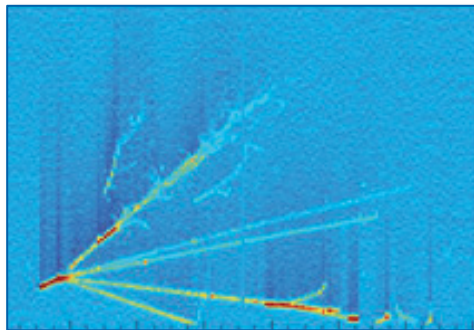


How are we different from CMS/CERN?

- Analysis Jobs go anywhere – the data follow
 - System does not pre-place nor pin data to sites
 - Utilizes wide variety of storage elements at remote sites
 - Fully automated (no human intervention)
- Fermilab tape archive (enstore) is the backend
 - Disk cache in front of tape system large enough to ensure long time to live for popular data
 - Obsolete data replaced by new/popular data automatically (LRU caching)
 - Fully automated (no human intervention)
 - Other tape systems (in2p3) and disk-only systems work also
- Operated by ~1/3 FTE
 - Rotating shift among data management group

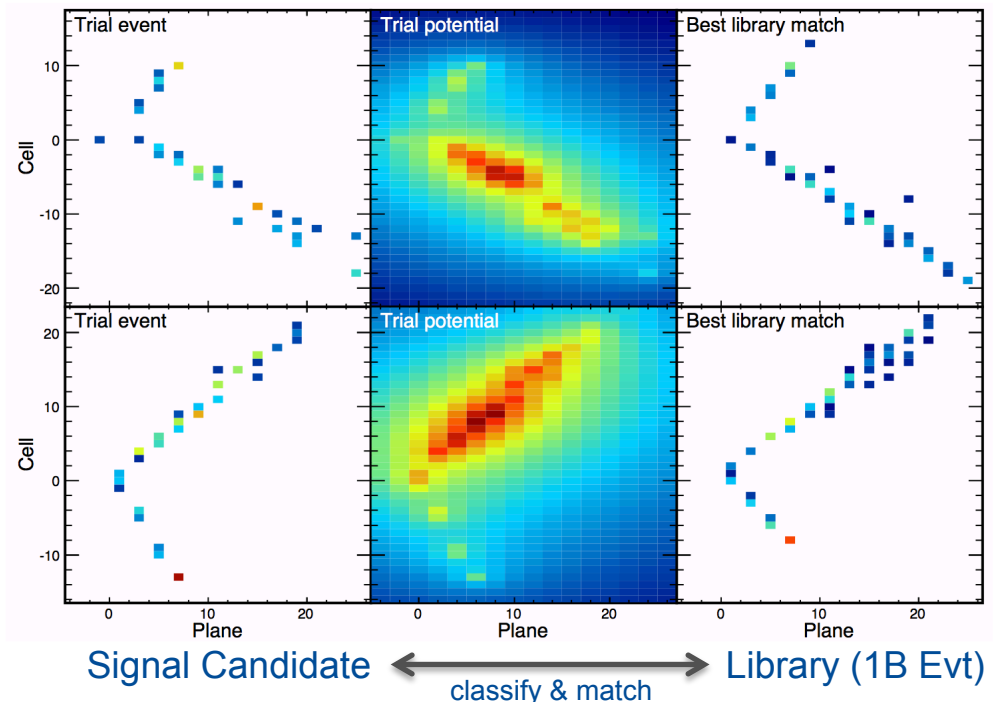
Beyond Big Data Research

- Our current solutions work well into the 100's PB/100's Million files range for “traditional” (embarrassingly parallel) analysis methods
- Next step is mapping new analysis methods to distributed data stores
 - High speed classification methods (template matching)
 - Map well onto map-reduce strategies



Liquid Ar Time Projection Chamber

Need to be able to perform
high resolution image processing
for next gen detectors (LArTPC)



Beyond Bigger Data

- Our current experiments (CMS, Atlas, Nova) write out event based data with 10's ns discretization.
- This works well for capturing “instantaneous” events but fails to capture long time scale correlations.
- Does the event model work for the next generation of experiments?
- Do we need to move to the “continuum limit” in how we treat our data to regain the long wavelength correlations?
- What about continuous readout technologies? How do we organize and analyze data in a grid/cluster environment without “chopping it up”?
- What does this mean for data storage systems?

Summary

- Particle Physics has a proven history of pushing the envelope of data analysis to make the next level of discovery
- We have a well defined solution to the current level of Big Data storage/analysis problems
 - Leverages metadata based catalogs with distributed object based data stores to provide high throughput data retrieval and analysis
 - Tools are developed, tested and proven to scale
- Methodology and tools map well into wider IoE domain of large distributed sensor networks
- Bigger Data analysis methods are under research for the physics domain as we move to real-time analysis of the continuum limit